

VIDEOITG: MULTIMODAL VIDEO UNDERSTANDING WITH INSTRUCTED TEMPORAL GROUNDING

Shihao Wang^{1*}, Guo Chen^{2*}, De-An Huang³, Zhiqi Li^{2*}, Minghan Li⁴, Guilin Liu³, Jose M. Alvarez³, Lei Zhang^{1†}, Zhiding Yu^{3†}

¹The Hong Kong Polytechnic Univ. ²Nanjing Univ. ³NVIDIA ⁴Harvard Univ.
<https://nvlabs.github.io/VideoITG/>

ABSTRACT

Recent studies have revealed that selecting informative and relevant video frames can significantly improve the performance of Video Large Language Models (Video-LLMs). Current methods, such as reducing inter-frame redundancy, employing separate models for image-text relevance assessment, or utilizing temporal video grounding for event localization, substantially adopt unsupervised learning paradigms, whereas they struggle to address the complex scenarios in long video understanding. We propose *Instructed Temporal Grounding for Videos (VideoITG)*, featuring customized frame sampling aligned with user instructions. The core of VideoITG is the *VidThinker* pipeline, an automated annotation framework that explicitly mimics the human annotation process. First, it generates detailed clip-level captions conditioned on the instruction; then, it retrieves relevant video segments through instruction-guided reasoning; finally, it performs fine-grained frame selection to pinpoint the most informative visual evidence. Leveraging VidThinker, we construct the VideoITG-40K dataset, containing 40K videos and 500K instructed temporal grounding annotations. We then design a plug-and-play VideoITG model, which takes advantage of visual language alignment and reasoning capabilities of Video-LLMs, for effective frame selection in a discriminative manner. Coupled with Video-LLMs, VideoITG achieves consistent performance improvements across multiple multimodal video understanding benchmarks, showing its superiority and great potentials for video understanding.

1 INTRODUCTION

The rapid advancement of Video Large Language Models (Video-LLMs) has opened new frontiers in video understanding, advancing complex tasks such as captioning [Chen et al. \(2024c\)](#); [Chai et al. \(2025\)](#); [Zhou et al. \(2024b\)](#); [Islam et al. \(2024\)](#); [Chen et al. \(2024d\)](#); [Wang et al. \(2024b\)](#), visual question answering [Fu et al. \(2024a\)](#); [Zhou et al. \(2024a\)](#); [Mangalam et al. \(2024\)](#); [Li et al. \(2024b\)](#); [Chen et al. \(2024a\)](#); [Xiao et al. \(2021\)](#); [Pătrăucean et al. \(2023\)](#), and embodied-agent applications [Brohan et al. \(2023\)](#); [Kim et al. \(2024\)](#); [Fu et al. \(2024b\)](#); [Liu et al. \(2024a\)](#); [Chen et al. \(2024b; 2025\)](#). However, these models face challenges when handling long videos due to high memory and computational demands. To mitigate this, existing approaches often adopt uniform frame sampling, a simple but naive strategy that frequently misses key frames for accurate video understanding, resulting in suboptimal performance.

To address these limitations, researchers have explored various strategies. One family of approaches focuses on reducing redundant spatiotemporal information by fusing or pruning overlapped content across frames, as can be seen in works employing pooling rules, similarity thresholds, or clustering to retain only essential frames [Xu et al. \(2024a\)](#); [Shen et al. \(2024\)](#); [Zhang et al. \(2024a\)](#); [Li et al. \(2024a\)](#); [Zhang et al. \(2024c\)](#). Another stream of strategies extends the length of model sequence to incorporate more tokens [Wang et al. \(2024c\)](#); [Team et al. \(2023\)](#), enabling longer temporal dependencies, despite the high computational cost, and risking information dilution. Alternative methods utilize question-

*Work done during an internship at NVIDIA.

†Corresponding authors.

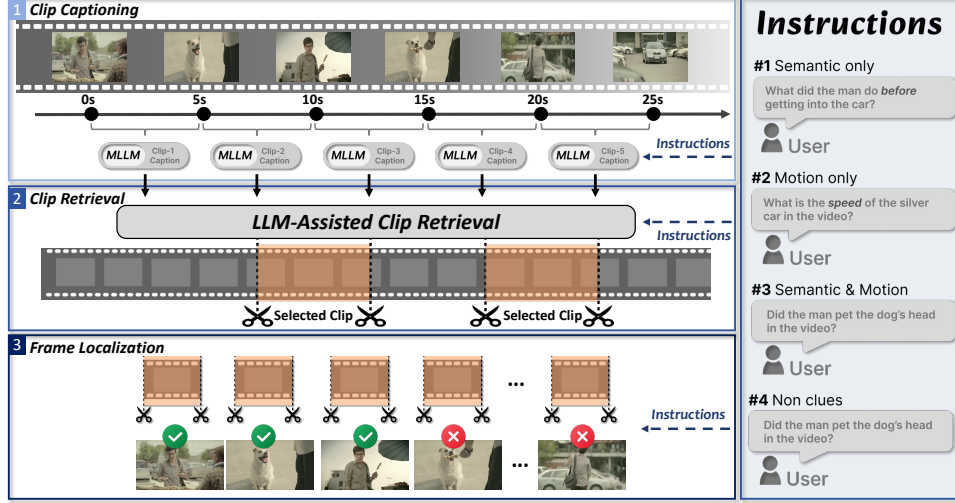


Figure 1: Overview of the VidThinker annotation pipeline for VideoITG. The pipeline consists of three stages that fully leverage the provided instructions: (1) segment-level clip captioning; (2) instruction-guided relevant clip retrieval; (3) fine-grained frame-level localization.

guided feature extraction or language queries to identify relevant segments [Li et al. \(2024c\)](#); [Yu et al. \(2023\)](#). For instance, SeViLA [Yu et al. \(2023\)](#) processes frames independently using BLIP-2 [Li et al. \(2023\)](#) before selecting keyframes, which serve as the input for subsequent video understanding tasks. However, the lack of temporal modeling capability limits its performance in tasks requiring multi-temporal cues.

Despite advances in compressing or extending the context for Video-LLMs, a performance gap persists between long and short videos due to limited training data for long-video content. When humans analyze long videos, they naturally employ a step-by-step approach: skimming the overall context, locating question-relevant clues, and then focusing on specific segments. Drawing inspiration from such a process, we propose **Instructed Temporal Grounding for Videos (VideoITG)**, which integrates user instructions into frame selection. While general temporal video grounding [Wang et al. \(2024a\)](#); [Qian et al. \(2024\)](#); [Lei et al. \(2021\)](#) emphasizes event localization within videos based on single temporal clue and descriptive language queries, VideoITG introduces a user-instruction-driven approach, customizing frame selection strategies to align with specific task requirements. Compared to existing frame selection frameworks [Yu et al. \(2023; 2025\)](#); [Han et al. \(2025\)](#), VideoITG can effectively handle multiple temporal clues for various tasks: localizing temporal cues from multiple clips to understand temporal relationships, employing event localization and uniform frame sampling to detect speed variations, and conducting diverse types of samplings to cover all videos for content captioning or existence judgment, *etc.*

To achieve the goal of VideoITG, we need to construct a comprehensive dataset for model training. To this end, we develop an automated data annotation pipeline, namely *VidThinker*, which consists of instruction-guided clip captioning, clip retrieval, and frame localization, ensuring high-quality task-aligned annotations, as shown in Fig. 1. Our annotation pipeline is inspired by the human reasoning process. It employs a coarse-to-fine strategy using GPT-4o [OpenAI \(2024\)](#) to generate detailed clip descriptions, followed by a “Needle-In-A-Haystack” approach for instruction-guided clip retrieval. To ensure a precise and instruction-aligned temporal grounding, we categorize instructions into four distinct types, each reflecting a unique reasoning requirement in video QA. Specifically, **Semantic-only** focuses on visual content such as objects or scenes to support appearance- or context-based questions. **Motion-only** targets dynamic cues such as movement or speed to answer motion-centric questions. **Semantic & Motion** integrates visual and temporal reasoning to handle questions involving appearance and motion jointly. **Non-clues** refers to general instructions without clear semantic or motion focus, aiming to maximize visual diversity for holistic understanding.

The resulting VideoITG-40K dataset contains 40K videos with varying durations (30s - 3mins) and 500K instruction-guided annotations, significantly surpassing previous temporal grounding datasets in both scale and quality of instruction-guided frame selection. Building on this foundation, we

present a family of VideoITG models that leverage text generation, anchor-based classification with causal attention, and pooling-based classification with full attention to enhance instructed temporal grounding and advance Video-LLM capabilities. In summary, our contributions are threefold:

- **VideoITG-40K dataset.** We define the tasks of VideoITG and develop an automated data annotation pipeline, namely *VidThinker*, to generate a large-scale dataset, namely VideoITG-40K, with 40K videos and 500K instruction-dependent annotations, allowing precise frame identification and effective video understanding.
- **VideoITG models.** We introduce a family of VideoITG models with varying attention and decoding strategies, designed to improve instruction-guided temporal grounding based on insights from the VideoITG-40K dataset.
- **Consistent improvement.** Our approach achieves consistent performance improvements on various multimodal video understanding benchmarks. By integrating VideoITG, we achieve improvements of **9.0%** on CG-Bench, **8.6%** on MLVU, **4.0%** on Video-MME, and **3.6%** on LongVideoBench for the InternVL2.5-8B model, showing the effectiveness of our framework.

2 RELATED WORK

2.1 VIDEO LARGE LANGUAGE MODELS

Recent advancements in video understanding with Video-LLMs introduce several strategies to address the temporal and spatial complexity associated with processing long videos. One key approach focuses on compressing visual features to make the process more efficient [Liu et al. \(2025\)](#); [Zohar et al. \(2024\)](#); [Ye et al. \(2024\)](#); [Wang et al. \(2024d\)](#). This can be done by using modules like Q-Former [Song et al. \(2024\)](#) and Perceiver Resampler [Zohar et al. \(2024\)](#), which merge frame features into a fixed number of queries. Spatial pooling techniques are also employed by models [Maaz et al. \(2024\)](#); [Xu et al. \(2024b;a\)](#) to effectively manage computational resources while preserving long temporal information. Another strategy aims to extend the effective sequence length to accommodate longer video inputs. For example, LongVA [Zhang et al. \(2024d\)](#), Qwen2-VL [Wang et al. \(2024c\)](#), and Gemini 1.5 Pro [Team et al. \(2023\)](#) increase token capacity, allowing for more extensive video analysis; however, this strategy often comes with higher computational costs [Wei & Chen \(2024\)](#); [Shu et al. \(2025\)](#) to manage the sequence length without auxiliary supervision.

Some methods focus on reducing interframe redundancy using similarity-based techniques, such as cosine similarity or clustering, to filter out redundant frames. Video-LaVIT [Jin et al. \(2024\)](#) and LongVU [Shen et al. \(2024\)](#) effectively apply these techniques to manage data more efficiently. However, many of these methods employ fixed thresholds or sampling strategies that might not be able to adequately capture the diversity and complexity inherent in real-world video content. In contrast, our proposed VideoITG offers a superior solution by utilizing instructed temporal grounding to align frame sampling with user instructions. This approach, combined with an automated annotation pipeline and a plug-and-play model, leads to significant performance improvements across multimodal video understanding benchmarks, demonstrating the effectiveness and scalability of VideoITG.

2.2 VIDEO TEMPORAL GROUNDING

Video Temporal Grounding [Ren et al. \(2024\)](#); [Wang et al. \(2024a\)](#); [Qian et al. \(2024\)](#); [Di & Xie \(2024\)](#) is a common task in video understanding that associates specific video moments with their corresponding timestamps, while Temporal Localization focuses on accurately identifying these moments within untrimmed videos [Liu et al. \(2024b\)](#); [Anne Hendricks et al. \(2017\)](#); [Li et al. \(2024d\)](#). Current Video-LLMs [Shen et al. \(2024\)](#); [Wang et al. \(2024a\)](#); [Huang et al. \(2024\)](#) have begun to leverage temporal grounding for frame selection by linking video content with temporal cues; however, existing methods [Huang et al. \(2025\)](#); [Yu et al. \(2023; 2025\)](#) mostly focus on single-time retrieval, which take descriptive annotations as input, limiting their generality and robustness in handling diverse real-world scenarios. Recognizing these limitations, we propose the VideoITG task, which introduces a customized sampling approach aligned with user instructions to improve the effectiveness of frame selection for a broad range of video understanding tasks.

3 VIDEOITG-40K: DATASET CONSTRUCTION

We first introduce *VidThinker* (Sec. 3.1), an automated annotation pipeline with three stages — clip captioning, retrieval, and frame localization — for instruction-based video annotation. We then detail our fine-grained instruction taxonomy and frame selection strategies (Sec. 3.2) to align annotations with QA task demands. Finally, we use *VidThinker* to construct the VideoITG-40K dataset and provide its statistics (Sec. 3.3).

3.1 VIDTHINKER: AUTOMATED ANNOTATION PIPELINE

When tackling instruction-driven temporal localization in long videos, human annotators typically follow a three-step reasoning process: they first parse the instruction to identify the question type and extract key information; then, they use these cues to narrow down the video to a coarse temporal window; finally, they perform fine-grained reasoning within that window to accurately pinpoint the occurrence of the target event. Inspired by this process, we propose *VidThinker*—an automated annotation pipeline designed to address the core challenge of instruction-guided temporal localization in long videos. *VidThinker* explicitly replicates the human step-by-step reasoning process, enabling fully automated, high-quality, and interpretable video annotations without manual labeling.

VidThinker decomposes the annotation process into three interdependent reasoning steps: i) Instructed Clip Captioning, ii) Instructed Clip Retrieval, and iii) Instructed Frame Localization. It progressively narrows the search space and enriches semantic alignment with the instruction.

i) Instructed Clip Captioning: The video v is uniformly divided into short clips (5 seconds each), denoted as $\{v_i\}_{i=0}^n$. For each segment, we employ LLM to extract salient phrases that capture the core information needed to fulfill the instruction. For example, given the question (q = ‘What does the man playing the drums do with his feet as he plays the drum?’) and the answer (a = ‘moves his feet’), the system distills the essential action phrase: k = ‘The man playing the drums moves his feet and hits the drums with his hands.’ We then input the extracted phrases alongside raw video clips into the MLLM to generate clip-level descriptions $\{c_i\}_{i=0}^n$ in a recurrent manner. The extracted phrases serve as reference cues to guide the model’s attention towards salient elements within each clip. However, the MLLM strictly adheres to visual evidence and it only incorporates information from the extracted phrases when it is explicitly observable in the current clip. This ensures that the system will not hallucinate or infer content solely based on the extracted phrases, maintaining descriptions grounded in visual content. The process can be formulated as follows:

$$k = \text{LLM}(q, a), \quad c_i = \text{MLLM}(k, v_i). \quad (1)$$

Conditioning on these instruction- and answer-derived cues, we ensure that the annotation of each segment is relevant and informative, thus facilitating precise instructed temporal grounding.

ii) Instructed Clip Retrieval: The generated clip descriptions $\{c_i\}_{i=0}^n$ are organized sequentially and evaluated by an LLM for the relevance to the QA pairs. Instead of simply assigning binary relevance scores, the LLM is instructed to perform chain-of-thought reasoning, explicitly considering both keyword matches and temporal relationships to directly output the indexes of relevant clips:

$$\mathcal{I}_{\text{rel-clip}} = \text{LLM}(\{c_i\}_{i=0}^n, q, a). \quad (2)$$

The chain-of-thought prompting requires the model to justify its selections based on both semantic and temporal cues, rather than relying solely on trivial keyword matching. This automation significantly improves the efficiency and the interpretability of relevant segment selection.

iii) Instructed Frame Localization: After coarse localization of video segment, *VidThinker* further refines the annotation by selecting key frames according to the instruction type. For each frame within the candidate segment, we prompt a large language model (LLM) to perform a binary classification task: given the QA pair and a single frame, the LLM determines whether the frame is relevant (yes) or not (no) to the instruction. Formally, for each frame f_i in the candidate segment, the LLM is prompted as follows:

$$y_i = \text{LLM}(f_i, q, a), \quad (3)$$

where $y_i \in \{\text{yes}, \text{no}\}$ indicates whether frame f_i is relevant to the QA pair. Only frames with positive responses ($y_i = \text{yes}$) are retained as the final temporal grounding results. This frame-level

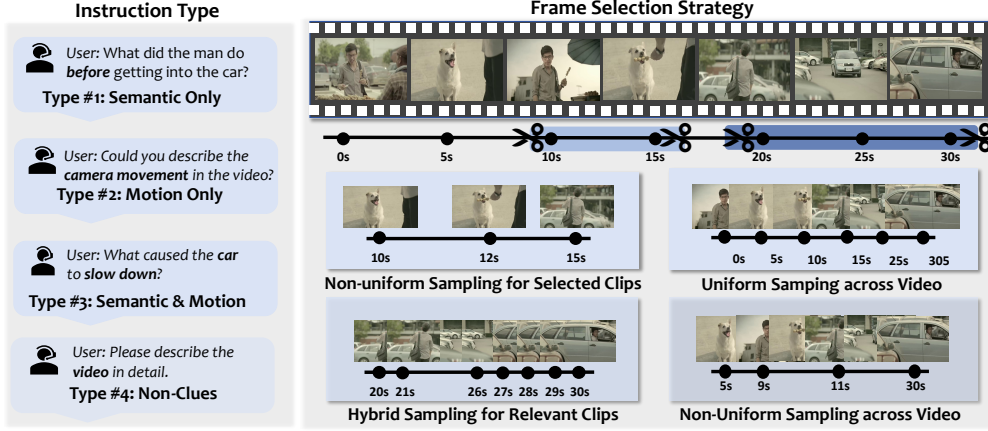


Figure 2: Illustration of four instruction types and their corresponding frame selection strategies in VidThinker. For semantic-focused instructions, the system selects diverse frames capturing key visual clues. For motion-focused instructions, frames are uniformly sampled to capture dynamic changes. When both semantic and motion cues are required, a hybrid sampling strategy is applied. For vague or open-ended instructions, the system samples a minimal yet diverse set of frames across the video for holistic coverage.

instruction-guided filtering allows *VidThinker* to achieve high precision in identifying the most informative frames for each instruction.

Leveraging the reasoning capabilities of MLLMs, *VidThinker* transforms video QA annotation into a fully automated, scalable, and cognitively inspired process. This approach not only reduces manual effort and variability, but also ensures high-quality, interpretable annotations suitable for training next-generation video understanding models.

3.2 FINE-GRAINED GROUNDING INSTRUCTION

To further enhance temporal grounding precision, we perform fine-grained frame selection tailored to different instructions types, ensuring that the final visual evidence aligns closely with the specific reasoning requirements of each QA task. Different types of instruction require different forms of visual understanding, such as static semantics, dynamic motion, or both—necessitating a customized frame selection strategies. Specifically, we categorize instructions into four types and apply corresponding frame sampling methods.

- **Semantic only:** These instructions focus on semantic content such as people, scenes, or objects. Following relevant segment localization, the system selects diverse frames that capture representative visual clues to ensure comprehensive semantic coverage. For example: “*What did the man do before getting into the car?*” The *VidThinker* needs to select frames that clearly show the man’s clothing and the guitar.
- **Motion only:** These instructions emphasize dynamic actions, such as movement type, speed, or direction. The frames are sampled at a fixed rate within the localized segment to accurately capture the progression of motion. For example: “*How does the person jump off the diving board?*” The *VidThinker* needs to select frames covering the sequence from takeoff, mid-air, to water entry.
- **Semantic & Motion:** These instructions require both semantic and motion understanding. The system applies fixed-rate sampling within motion-relevant regions while ensuring the preservation of semantically informative frames, balancing both needs. For example: “*Could you describe the camera movement in the video?*” The *VidThinker* needs to select frames showing hand drumming and foot movement simultaneously.
- **Non Clues:** These instructions are open-ended or vague without clear semantic or motion focus, aiming to maximize visual diversity for holistic understanding. In these cases, the system selects a small yet diverse set of frames across the entire video to maximize visual information

Table 1: Comparison of dataset statistics for temporal grounding and highlight detection datasets.

Dataset	# Videos	# Queries	Avg. Duration	Instructed?
DiDeMo Anne Hendricks et al. (2017)	10.6K	41.2K	29s	No
QuerYD Oncescu et al. (2021)	2.6K	32K	278s	No
HiREST Zala et al. (2023)	3.4K	8.6K	263s	No
Charades-STA Gao et al. (2017)	6.7K	16.1K	30s	No
QVHighlights Lei et al. (2021)	10.2K	10.3K	150s	No
VideoITG-40K	40K	500K	120s	Yes

coverage while minimizing redundancy. For example: “Please describe the video in detail.” The *VidThinker* selects representative frames from the beginning, middle, and end of the video.

3.3 DATASET STATISTICS

Leveraging our proposed *VidThinker* pipeline, we construct the VideoITG-40K dataset, which is sourced from the LLaVA-Video dataset [Zhang et al. \(2024e\)](#). VideoITG-40K achieves an unprecedented scale, comprising 40,000 videos and 500,000 annotations tailored specifically for instruction-guided temporal grounding. The entire annotation process is automatically carried out by *VidThinker*, ensuring high efficiency, consistency, and alignment with diverse instruction types.

VideoITG-40K contains videos of varying duration, averaging 120 seconds, and is uniformly sampled across the timelength of 30-60s, 1-2mins, and 2-3mins. Each video is comprehensively annotated with 10–15 QA pairs, including both multiple-choice and open-ended questions. As shown in Table 1, VideoITG-40K significantly surpasses existing datasets in volume, with nearly four times the number of videos compared to DiDeMo [Anne Hendricks et al. \(2017\)](#) (10.6K) and QVHighlights [Lei et al. \(2021\)](#) (10.2K), and far exceeding others like QuerYD [Oncescu et al. \(2021\)](#) (2.6K) and HiREST [Zala et al. \(2023\)](#) (3.4K). Unlike prior datasets that primarily focus on descriptive text queries for video understanding, VideoITG-40K distinguishes itself through its instruction-guided approach, enabling models to locate relevant video content based on specific user queries.

4 VIDEOITG: MODEL DESIGN

In this section, we explore how to utilize our VideoITG-40K dataset to train the model for the **Instructed Temporal Grounding** task, aiming to optimize video frame selection and enhance the performance of Video-LLMs. Our framework, as shown in Fig. 3, consists of three main components: a vision encoder (*i.e.*, VIT) for extracting text-aligned visual features F , a VideoITG model for instruction-guided frame selection \mathcal{I}_{rel} , and a VideoLLM for generating answers a based on the selected frames $F_{\mathcal{I}_{\text{rel}}}$ and the question q . The process can be described as follows:

$$F = \text{VIT}(v), \quad (4)$$

$$\mathcal{I}_{\text{rel}} = \text{VideoITG}(F, q), \quad (5)$$

$$a = \text{VideoLLM}(F_{\mathcal{I}_{\text{rel}}}, q). \quad (6)$$

The VideoITG model is designed in a plug-and-play fashion, which is focused on the following aspects: (1) to which extent the Video-LLMs capitalize on the alignment between video and language tokens, as well as their ability to follow instructions; (2) whether the model has sufficient contextual encoding capabilities to handle and analyze multiple temporal cues; and (3) the model’s capacity to manage both single-turn and multi-turn conversations. With the above considerations, we develop three model variants: text generation-based classification, anchor-based classification, and pooling-based classification, as illustrated in Fig. 3 (b).

Variant A: text-generation-based classification. As shown on the left side of Fig. 3 (b), we start by discussing the text generation design, where Instructed Temporal Grounding is framed as a next-token prediction task, producing text tokens as output. This approach aligns with the current training paradigm of Video-LLMs, optimizing the use of Video-LLM for vision-language alignment and instruction following. It also supports multi-turn dialogue without repeatedly encoding the visual features. Previous works, such as Timechat [Ren et al. \(2024\)](#) and Grounded-VideoLLM [Wang et al. \(2024a\)](#), use this paradigm to tackle time-sensitive tasks. However, this method has some drawbacks,

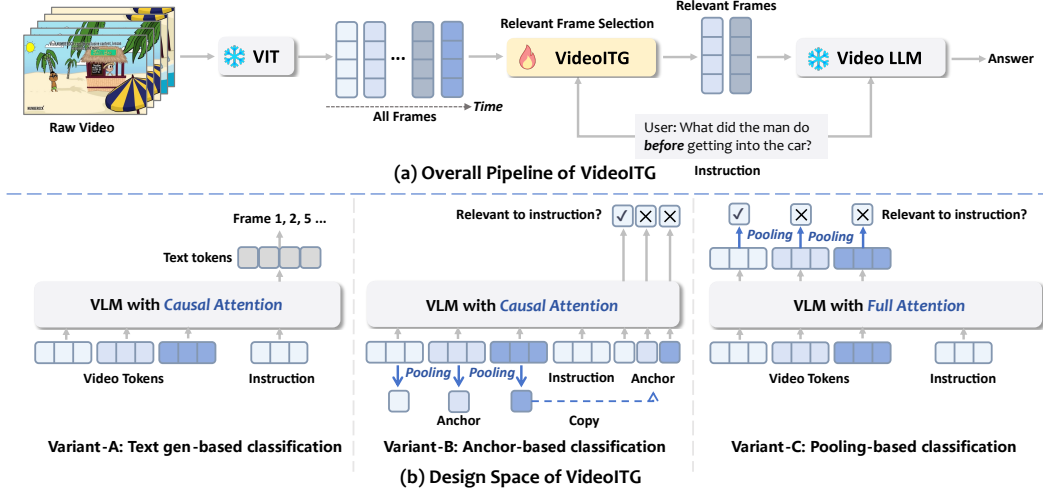


Figure 3: **VideoITG model design:** (a) Text generation aligns video and language tokens for sequential predictions. (b) Classification with causal attention utilizes anchor tokens for temporal cue management. (c) Classification with full attention facilitates interaction across visual and text tokens without anchors.

such as the inefficiency during inference and the shortcut learning due to the teacher forcing strategy during training.

Variant B: anchor-based classification. Another model design adopts a discriminative paradigm by classifying the visual tokens corresponding to each video frame, as shown in the middle of Fig. 3 (b). We initialize the model with a Video-LLM and retain the causal attention mask to remain consistent with the original Video-LLM design, thereby preserving its pre-training capabilities. However, the causal attention mask prevents visual tokens from accessing the instruction in advance, and earlier frame features are unable to access subsequent frame features, constraining the model’s ability to handle multiple temporal cues. To address this limitation, we introduce an *anchor token* after the instruction. For a video frame at time index t , we compute an anchor token A^t via global average pooling over all spatial positions, as formulated below:

$$A^t = \frac{1}{M} \sum_{i,j} F_{ij}^t, \quad (7)$$

where F_{ij}^t represents the extracted visual feature at 2D grid position (i, j) of the t -th frame, and M is the total number of patches within each frame. For a video with T frames, we compute T anchor tokens in total: $\{A^t\}_{t=1}^T$. This paradigm is friendly to multi-turn conversations because retaining the causal attention mask structure permits the efficient use of KV-cache for optimization.

Variant C: pooling-based classification. As discussed above, the presence of a causal attention mask makes it difficult to directly supervise the classification of visual tokens for each frame. Therefore, we also attempt to remove the causal attention mask, allowing visual tokens and text instruction tokens to interact through full attention across the sequence. Following this, we perform average pooling and classification on the visual tokens of each frame without establishing separate anchor tokens. The overall process is illustrated in the right of Fig. 3 (b). The main drawback of this paradigm is that in multi-turn conversations scenarios, the model cannot reuse visual tokens through KV-Cache.

5 EXPERIMENTS

5.1 IMPLEMENTATION DETAILS

We follow the training approach of LLaVA-Video [Zhang et al. \(2024e\)](#), using the pretrained model as the initialization for our VideoITG model’s pre-training. We employ SigLIP [Zhao et al. \(2023\)](#) as the vision encoder and Qwen2 [Wang et al. \(2024c\)](#) as the language model. Initially, we train the

Table 2: Performance comparison of VideoITG integrated with different Video-LLMs, varying in both the size of the answering LLM and the number of sampled frames. “Uni- k ” denotes uniform sampling of k frames, while “Top- k ” refers to selecting the top k frames based on relevance scores generated by our proposed VideoITG.

LLM	Selection	LongVideoBench	MLVU	VideoMME			CG-Bench	Avg.
		8min	12min	S (2 min)	M (10 min)	L (40 min)	27min	
LLaVA-Video-7B	Uni-32	58.7	66.8	76.3	60.3	52.7	35.8	58.4
	Top-32 (Ours)	61.6	74.6	77.3	65.9	55.2	42.8	62.9
	Δ	+2.9%	+7.8%	+1.0%	+5.6%	+2.5%	+7.0%	+4.5%
LLaVA-Video-7B	Uni-64	59.9	70.2	75.8	63.0	54.7	36.9	60.1
	Top-64 (Ours)	60.9	76.3	76.1	66.0	57.0	42.9	63.2
	Δ	+1.0%	6.1%	+0.3%	+3.0%	+2.3%	+6.0%	+3.1%
InternVL2.5-8B	Uni-32	58.3	66.4	75.1	61.7	53.1	37.7	58.7
	Top-32 (Ours)	61.9	75.0	78.0	67.1	56.9	46.7	64.3
	Δ	+3.6%	+8.6%	+2.9%	+5.4%	+3.8%	+9.0%	+5.6%
InternVL2.5-26B	Uni-32	55.6	71.3	78.1	67.1	56.9	40.6	61.6
	Top-32 (Ours)	63.0	78.9	80.8	69.0	59.9	48.7	66.7
	Δ	+7.4%	+7.6%	+2.7%	+1.9%	+3.0%	+8.1%	+5.1%

MLP projector on image caption datasets with a batch size of 256 and a learning rate of 1×10^{-3} . Then, we fine-tune all model parameters on the LLaVA-OV-SI Li et al. (2024a) and LLaVA-Video datasets. During this stage, the video frame sampling rate is set to 64, and the LLM’s maximum sequence length is set to 16K. We then train the VideoITG model on the proposed VideoITG-40K dataset, adjusting the video sampling rate to 1 fps.

Throughout training and inference, we employ a dynamic token spatial size strategy Liu et al. (2025). Across all stages, the LLM’s learning rate is 2×10^{-5} , and in the final stage, the learning rate for the classification head is 2×10^{-4} . To fairly compare with other leading video LLMs, we primarily use results from their original papers. When results are unavailable, we integrate the models into LLMs-Eval Zhang et al. (2024b) and assess them under consistent settings. Due to context length constraints, we support up to 512 video frames as input (with 16 visual tokens per frame) for the VideoITG model, from which we select the top 32 frames based on their scores by default. We evaluate on three long video datasets: LongVideoBench Wu et al. (2024), MLVU Zhou et al. (2024a), and VideoMME Fu et al. (2024a). For more results, please see the Appendix.

5.2 MAIN RESULTS

In Table 2, we integrate our VideoITG model with various Video-LLMs to examine how different frame sampling strategies and the number of sampled frames influence answer quality. As can be seen, VideoITG’s frame selection strategy significantly outperforms uniform sampling across both 32-frame and 64-frame settings. This demonstrates that uniform sampling indeed constrains the extraction of informative content within the limited frame budget.

Furthermore, we evaluate the Video-LLM models across different sizes. For the InternVL2.5-26B model, despite having a higher baseline compared to InternVL2.5-8B, it still demonstrates substantial performance improvements, with a 7.4% increase on LongVideoBench, a 7.6% increase on MLVU and a 9.0 % improvement on CG-Bench, highlighting the effectiveness of our approach even with more advanced models. Notably, when integrating VideoITG with InternVL2.5-8B, the performance reaches 64.3% on average, which surpasses the baseline performance of InternVL2.5-26B (61.6%). This is particularly significant as it shows that a smaller model with intelligent frame selection can outperform a much larger model using standard uniform sampling. The performance gains are especially pronounced on longer video benchmarks, where the InternVL2.5-8B model with VideoITG achieves 46.7% on CG-Bench, exceeding the InternVL2.5-26B baseline of 40.6%. These findings suggest that effective frame selection can be more impactful than simply scaling up model size, particularly for long-video understanding tasks.

In Table 3, we compare the performance of different models on various video benchmarks. We report the results of InternVL2.5-8B integrated with VideoITG. Our InternVL2.5-8B-ITG model achieves state-of-the-art performance across multiple datasets, demonstrating the effectiveness of our proposed

Table 3: The performance (accuracy) of SOTA methods on video benchmarks. For InternVL2.5-8B results, we report the higher results in the technical report and Imms-eval. We sample 32 frames using VideoITG for our results.

Model	Open-Ended Q&A		Multi-Choice Q&A						
	ActNet-QA	EgoSchema	MLVU	NExT-QA	Perception Test	LongVideoBench	VideoMME	MLVBench	
	test	test	m-avg	mc	val	val	wo/w-sub	val	
<i>Proprietary models</i>									
GPT-4o (OpenAI, 2024)	-	-	64.6	-	-	66.7	71.9/77.2	43.5	
Gemini-1.5-Flash (Team et al., 2023)	55.3	65.7	-	-	-	61.6	70.3/75.0	-	
Gemini-1.5-Pro (Team et al., 2023)	57.5	72.2	-	-	-	64.0	75.0/81.3	-	
<i>Open-source models</i>									
VILA-40B (Lin et al., 2024)	58.0	58.0	-	67.9	54.0	-	60.1/61.1	-	
PLLaVA-34B (Xu et al., 2024a)	60.9	-	-	-	-	53.2	-	58.1	
VideoLLaMA2-7B (Cheng et al., 2024)	50.2	50.5	-	-	49.6	-	45.1/46.6	53.4	
LongVA-7B (Zhang et al., 2024d)	50.0	-	56.3	68.3	-	-	52.6/54.3	-	
LongVU-7B (Zhang et al., 2024d)	-	67.6	65.4	-	-	-	60.6/-	66.9	
LLaVA-OV-7B (Li et al., 2024a)	56.6	60.1	64.7	79.4	57.1	56.5	58.2/61.5	56.7	
mPLUG-Owl3-8B (Ye et al., 2024)	-	-	-	78.6	-	52.1	53.5/-	54.5	
LLaVA-Video-7B (Zhang et al., 2024e)	56.5	57.3	70.8	83.2	67.9	58.2	63.3/69.7	58.6	
Qwen2-VL-7B (Wang et al., 2024c)	-	66.7	-	-	62.3	55.6	63.3/69.0	67.0	
InternVL2.5-8B (Zhang et al., 2024c)	-	51.5	68.9	-	-	60.0	64.2/66.9	72.0	
InternVL2.5-8B-ITG	57.4	51.6	75.0	79.5	64.9	61.9	67.3/69.6	72.2	

VideoITG framework. Particularly, in long-video understanding tasks, our model exhibits significant advantages, outperforming current SOTA models like Qwen2-VL by substantial margins.

5.3 ABLATION ON VIDEOITG DESIGN CHOICES

Table 4 presents a comprehensive analysis on the design of our VideoITG framework, directly supporting our key contributions.

Architecture. First, we compare the three variants of our model architecture. We observe that Variant A, which is based on the text generation paradigm, performs the worst. One possible reason is that text generation models trained with the next-token prediction paradigm suffer from sparse supervision due to teacher forcing, where previous frame selections influence subsequent ones, making the training process less efficient compared to discriminative classification models. We find that Variant C with full-attention outperforms Variant B with causal attention. This improvement may be attributed to full-attention’s larger receptive field, which enables global temporal relationship modeling and allows all tokens to access the textual query simultaneously.

Dataset. We analyze our data annotation strategies to demonstrate the effectiveness of our *Vid-Thinker* annotation pipeline. Ablation studies show that the performance degrades when Instructed Clip Captioning are removed, with accuracy dropping from 56.9% to 53.4% on Videomme Long videos and from 75.0% to 73.2% on MLVU. This demonstrates that ensuring information diversity is crucial for maintaining comprehensive feature representation of videos. Similarly, removing Instructed Frame Localization decreases performance, particularly on Videomme Medium videos (from 67.1% to 65.8%). These results confirm that both stages are essential for optimal model performance and validate our data construction approach of the VideoITG-40K dataset.

Pre-training. Finally, we investigate the impact of vision-language alignment pre-training on model performance. Our experiments reveal that removing video pre-training causes only modest performance changes across benchmarks, with slight increases on Videomme Long videos. This suggests that the benefits of video data for instructed temporal grounding tasks primarily stem from effective visual context length, yet this impact is relatively minor compared to vision-language alignment. This observation is further validated if we eliminate both image and video pre-training data, starting from a text-only large language model, where performance drops dramatically, with accuracy decreasing from 75.0% to 69.1% on MLVU and from 61.9% to 58.6% on LongVideoBench.

Table 4: Empirical studies on the VideoITG-40k dataset and VideoITG model design. We adopt Variant-C for subsequent experiments. “No Images” and “No Videos” indicate that image-text data (LAION-CC-SBU-558K & LLaVA-OV-SI) or video data (LLaVA-Video-178K) are excluded from pre-training, respectively.

Abaltion	Experiment	Videomme			MLVU (%) ↑	LongVideoBench (%) ↑
		Short (%) ↑	Medium (%) ↑	Long (%) ↑		
Architecture	Variant-A-7B	51.0	44.8	44.4	45.7	56.8
	Variant-B-7B	77.9	66.0	56.2	74.6	61.3
	Variant-C-7B	78.0	67.1	56.9	75.0	61.9
	Variant-C-3B	77.1	64.8	56.0	74.5	61.5
Dataset Construction	No Clip Captioning	77.5	63.1	53.4	73.2	61.7
	No Frame Localization	77.6	65.8	56.8	74.1	61.5
Pre-training Data	No Videos	77.2	64.9	57.4	74.5	61.6
	No Images & Videos	76.6	63.0	54.4	69.1	58.6

Table 5: Results with different scoring LMMs. SigLIP performs worse on three out of four benchmarks, with an average score similar to uniform sampling. In the third row, we follow [Huang et al. \(2025\)](#) and use a standalone VLM to assess the relevance between the question and each frame, and select the top 32 frames with the highest probability of “Yes” as the output.

Selection Methods	Answering LMM	Frames	LongVideoBench	MLVU	VideoMME	Avg.
None (Uniform)	InternVL2.5-8B	32	58.3	66.4	63.3	62.7
SigLIP Zhai et al. (2023)	InternVL2.5-8B	32	60.4	69.3	62.4	64.0
InternVL2.5-8B Huang et al. (2025)	InternVL2.5-8B	32	60.7	70.3	64.7	65.2
VideoITG-7B (Ours)	InternVL2.5-8B	32	61.9	75.0	67.3	68.1

This substantial degradation underscores that robust vision-language alignment is crucial to effective VideoITG training.

5.4 ABLATION ON SELECTION METHODS

We compare different kinds of frame selection methods in Table 5. The uniform sampling method, serving as a baseline, achieves an average score of 62.7. SigLIP slightly improves this score to 64.0. However, SigLIP is limited to handling shorter descriptive text and lacks the capability to understand complex user instructions, resulting in poorer performance. When using InternVL2.5-8B as both the selection and answering LMM [Huang et al. \(2025\)](#), there is a further improvement, yielding an average score of 65.2. This approach leverages a standalone VLM to assess the relevance between each frame and the given question. While this method benefits from the VLM’s knowledge and instruction-following abilities, it falls short in temporal modeling and performs poorly when multiple temporal cues are involved. Moreover, processing each frame independently is inefficient. Notably, the VideoITG-7B model outperforms all other methods with an average score of 68.1. Our proposed VideoITG framework excels because both data annotation and model design are aware of instruction following and temporal modeling.

5.5 VISUALIZATION

In Fig. 4, we present two cases comparing uniform sampling versus VideoITG sampling of 8 frames from the VideoMME [Fu et al. \(2024a\)](#) Benchmark. The first case demonstrates a multi-temporal cues problem requiring temporal relationship understanding between brushing teeth and spraying perfume actions. VideoITG effectively captures both actions, enabling correct temporal ordering determination. The second case focuses on identifying the final action in the video, where VideoITG precisely captures the rapid consecutive movements at the end of the sequence. In both cases, uniform sampling either misses critical temporal relationships or fails to capture rapid action sequences, leading to incomplete or incorrect video understanding.

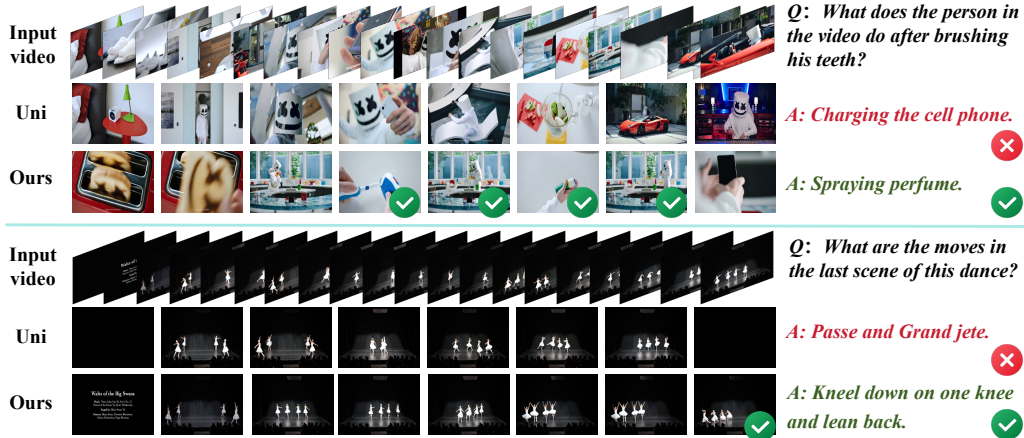


Figure 4: Two examples of how different sampling strategies impact video understanding. We mark the identified key frames that directly answer the question with green check-marks.

6 CONCLUSION

In this paper, we presented VideoITG, a novel framework for instruction-aligned frame selection in Video-LLMs. The key to our approach was the *VidThinker* pipeline, which mimics human annotation by generating detailed, instruction-guided clip descriptions, retrieving relevant segments, and performing fine-grained frame selection. Using this pipeline, we constructed the VideoITG-40K dataset with 40K videos and 500K temporal grounding annotations. Based on this resource, we developed plug-and-play VideoITG models that leverage visual-language alignment and reasoning to handle diverse temporal grounding tasks. Experiments showed that VideoITG consistently improves Video-LLMs’ performance across multiple video understanding benchmarks, highlighting its effectiveness and potential for advancing instruction-driven video understanding.

Limitations. Our current framework consists of two separate modules during inference: VideoITG for frame selection and a standalone Video-LLM for question answering. Although we have carefully designed the frame labeling strategies, the lack of gradient optimization between these two modules during training can lead to suboptimal results. Our VideoITG framework serves as a promising starting point, while in future work we could explore reinforcement learning techniques to bridge these two modules, enabling more efficient and accurate frame selection.

REFERENCES

- Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language. In *ICCV*, 2017. 3, 6
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. RT-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv:2307.15818*, 2023. 1
- Wenhao Chai, Enxin Song, Yilun Du, Chenlin Meng, Vashisht Madhavan, Omer Bar-Tal, Jeng-Neng Hwang, Saining Xie, and Christopher D Manning. AuroraCap: Efficient, performant video detailed captioning and a new benchmark. 2025. 1
- Guo Chen, Yicheng Liu, Yifei Huang, Yuping He, Baoqi Pei, Jilan Xu, Yali Wang, Tong Lu, and Limin Wang. CG-Bench: Clue-grounded question answering benchmark for long video understanding. *arXiv:2412.12075*, 2024a. 1
- Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong Lin, Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng Gao, Dongxing Mao, and Mike Zheng Shou. VideoLLM-online: Online video large language model for streaming video. In *CVPR*, 2024b. 1

-
- Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Bin Lin, Zhenyu Tang, et al. ShareGPT4Video: Improving video understanding and generation with better captions. *arXiv:2406.04325*, 2024c. 1
- Qirui Chen, Shangzhe Di, and Weidi Xie. Grounded multi-hop videoqa in long-form egocentric videos. In *AAAI*, 2025. 1
- Yang Chen, Sheng Guo, and Limin Wang. A large-scale study on video action dataset condensation. *arXiv:2412.21197*, 2024d. 1
- Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. VideoLLaMA 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv:2406.07476*, 2024. 9
- Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *ICLR*, 2024. 16
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *NeurIPS*, 2022. 16
- Shangzhe Di and Weidi Xie. Grounded question-answering in long egocentric videos. In *CVPR*, 2024. 3
- Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-MME: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv:2405.21075*, 2024a. 1, 8, 10
- Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Shaoqi Dong, Xiong Wang, Di Yin, Long Ma, et al. VITA: Towards open-source interactive omni multimodal llm. *arXiv:2408.05211*, 2024b. 1
- Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. TALL: Temporal activity localization via language query. In *ICCV*, 2017. 6
- Songhao Han, Wei Huang, Hairong Shi, Le Zhuo, Xiu Su, Shifeng Zhang, Xu Zhou, Xiaojuan Qi, Yue Liao, and Si Liu. Videospresso: A large-scale chain-of-thought dataset for fine-grained video reasoning via core frame selection. In *CVPR*, 2025. 2
- De-An Huang, Shijia Liao, Subhashree Radhakrishnan, Hongxu Yin, Pavlo Molchanov, Zhiding Yu, and Jan Kautz. LITA: Language instructed temporal-localization assistant. In *ECCV*, 2024. 3
- De-An Huang, Subhashree Radhakrishnan, Zhiding Yu, and Jan Kautz. FRAG: Frame selection augmented generation for long video and long document understanding. *arXiv:2504.17447*, 2025. 3, 10
- Md Mohaiminul Islam, Ngan Ho, Xitong Yang, Tushar Nagarajan, Lorenzo Torresani, and Gedas Bertasius. Video Recap: Recursive captioning of hour-long videos. In *CVPR*, 2024. 1
- Yang Jin, Zhicheng Sun, Kun Xu, Liwei Chen, Hao Jiang, Quzhe Huang, Chengru Song, Yuliang Liu, Di Zhang, Yang Song, et al. Video-LaVIT: Unified video-language pre-training with decoupled visual-motional tokenization. *arXiv:2402.03161*, 2024. 3
- Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. OpenVLA: An open-source vision-language-action model. *arXiv:2406.09246*, 2024. 1
- Jie Lei, Tamara L Berg, and Mohit Bansal. Detecting moments and highlights in videos via natural language queries. In *NeurIPS*, 2021. 2, 6
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. LLaVA-OneVision: Easy visual task transfer. *arXiv:2408.03326*, 2024a. 1, 8, 9
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 2

-
- Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. MVBench: A comprehensive multi-modal video understanding benchmark. In *CVPR*, 2024b. 1
- Yanwei Li, Chengyao Wang, and Jiaya Jia. LLaMA-VID: An image is worth 2 tokens in large language models. In *ECCV*, 2024c. 2
- Zeqian Li, Qirui Chen, Tengda Han, Ya Zhang, Yanfeng Wang, and Weidi Xie. Multi-sentence grounding for long-term instructional video. In *ECCV*, 2024d. 3
- Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. VILA: On pre-training for visual language models. In *CVPR*, 2024. 9
- Jihao Liu, Zhiding Yu, Shiyi Lan, Shihao Wang, Rongyao Fang, Jan Kautz, Hongsheng Li, and Jose M Alvarez. StreamChat: Chatting with streaming video. *arXiv:2412.08646*, 2024a. 1
- Ye Liu, Zongyang Ma, Zhongang Qi, Yang Wu, Ying Shan, and Chang Wen Chen. E.T. Bench: Towards open-ended event-level video-language understanding. In *NeurIPS*, 2024b. 3
- Zuyan Liu, Yuhao Dong, Ziwei Liu, Winston Hu, Jiwen Lu, and Yongming Rao. Oryx MLLM: On-demand spatial-temporal understanding at arbitrary resolution. In *ICLR*, 2025. 3, 8
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-ChatGPT: Towards detailed video understanding via large vision and language models. In *ACL*, 2024. 3
- Kartikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. EgoSchema: A diagnostic benchmark for very long-form video language understanding. In *NeurIPS*, 2024. 1
- Andreea-Maria Oncescu, Joao F Henriques, Yang Liu, Andrew Zisserman, and Samuel Albanie. QuerYD: A video dataset with high-quality text and audio narrations. In *ICASSP*, 2021. 6
- OpenAI. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o/>, 2024. 2, 9, 16
- Viorica Pătrăucean, Lucas Smaira, Ankush Gupta, Adrià Recasens Contente, Larisa Markeeva, Dylan Banarse, Skanda Koppula, Joseph Heyward, Mateusz Malinowski, Yi Yang, Carl Doersch, Tatiana Matejovicova, Yury Sulsky, Antoine Miech, Alex Frechette, Hanna Klimczak, Raphael Koster, Junlin Zhang, Stephanie Winkler, Yusuf Aytar, Simon Osindero, Dima Damen, Andrew Zisserman, and João Carreira. Perception Test: A diagnostic benchmark for multimodal video models. In *NeurIPS*, 2023. 1
- Long Qian, Juncheng Li, Yu Wu, Yaobo Ye, Hao Fei, Tat-Seng Chua, Yueting Zhuang, and Siliang Tang. Momentor: Advancing video large language model with fine-grained temporal reasoning. In *ICML*, 2024. 2, 3
- Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu Hou. TimeChat: A time-sensitive multimodal large language model for long video understanding. In *CVPR*, 2024. 3, 6
- Xiaoqian Shen, Yunyang Xiong, Changsheng Zhao, Lemeng Wu, Jun Chen, Chenchen Zhu, Zechun Liu, Fanyi Xiao, Balakrishnan Varadarajan, Florian Bordes, Zhuang Liu, Hu Xu, Hyunwoo J. Kim, Bilge Soran, Raghuraman Krishnamoorthi, Mohamed Elhoseiny, and Vikas Chandra. LongVU: Spatiotemporal adaptive compression for long video-language understanding. *arXiv:2410.17434*, 2024. 1, 3
- Yan Shu, Zheng Liu, Peitian Zhang, Minghao Qin, Junjie Zhou, Zhengyang Liang, Tiejun Huang, and Bo Zhao. Video-XL: Extra-long vision language model for hour-scale video understanding. In *CVPR*, 2025. 3
- Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. MovieChat: From dense token to sparse memory for long video understanding. In *CVPR*, 2024. 3
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv:2312.11805*, 2023. 1, 3, 9

-
- Haibo Wang, Zhiyang Xu, Yu Cheng, Shizhe Diao, Yufan Zhou, Yixin Cao, Qifan Wang, Weifeng Ge, and Lifu Huang. Grounded-VideoLLM: Sharpening fine-grained temporal grounding in video large language models. *arXiv:2410.03290*, 2024a. 2, 3, 6
- Jiawei Wang, Liping Yuan, Yuchen Zhang, and Haomiao Sun. Tarsier: Recipes for training and evaluating large video description models. *arXiv:2407.00634*, 2024b. 1
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-VL: Enhancing vision-language model’s perception of the world at any resolution. *arXiv:2409.12191*, 2024c. 1, 3, 7, 9
- Yuxuan Wang, Cihang Xie, Yang Liu, and Zilong Zheng. VideoLLaMB: Long-context video understanding with recurrent memory bridges. *arXiv:2409.01071*, 2024d. 3
- Hongchen Wei and Zhenzhong Chen. Visual context window extension: A new perspective for long video understanding. *arXiv:2409.20018*, 2024. 3
- Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. LongVideoBench: A benchmark for long-context interleaved video-language understanding. In *NeurIPS*, 2024. 8
- Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. NExT-QA: Next phase of question-answering to explaining temporal actions. In *CVPR*, 2021. 1
- Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See Kiong Ng, and Jiashi Feng. PLLaVA: Parameter-free llava extension from images to videos for video dense captioning. *arXiv:2404.16994*, 2024a. 1, 3, 9
- Mingze Xu, Mingfei Gao, Zhe Gan, Hong-You Chen, Zhengfeng Lai, Haiming Gang, Kai Kang, and Afshin Dehghan. SlowFast-LLaVA: A strong training-free baseline for video large language models. *arXiv:2407.15841*, 2024b. 3
- Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mPLUG-Owl3: Towards long image-sequence understanding in multi-modal large language models. In *ICLR*, 2024. 3, 9
- Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. Self-chained image-language model for video localization and question answering. In *NeurIPS*, 2023. 2, 3
- Sicheng Yu, Chengkai Jin, Huanyu Wang, Zhenghao Chen, Sheng Jin, Zhongrong Zuo, Xiaolei Xu, Zhenbang Sun, Bingni Zhang, Jiawei Wu, et al. Frame-Voyager: Learning to query frames for video large language models. In *ICLR*, 2025. 2, 3
- Abhay Zala, Jaemin Cho, Satwik Kottur, Xilun Chen, Barlas Oguz, Yashar Mehdad, and Mohit Bansal. Hierarchical video-moment retrieval and step-captioning. In *CVPR*, 2023. 6
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *ICCV*, 2023. 7, 10
- Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi Feng, Jifeng Dai, and Xiaojie Jin. Flash-VStream: Memory-based real-time understanding for long video streams. *arXiv:2406.08085*, 2024a. 1
- Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, et al. LMMs-Eval: Reality check on the evaluation of large multimodal models. *arXiv:2407.12772*, 2024b. 8
- Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong Duan, Bin Wang, Linke Ouyang, et al. InternLM-XComposer-2.5: A versatile large vision language model supporting long-contextual input and output. *arXiv:2407.03320*, 2024c. 1, 9
- Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to vision. *arXiv:2406.16852*, 2024d. 3, 9

-
- Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv:2410.02713*, 2024e. [6](#), [7](#), [9](#), [16](#)
- Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. MLVU: A comprehensive benchmark for multi-task long video understanding. *arXiv:2406.04264*, 2024a. [1](#), [8](#)
- Xingyi Zhou, Anurag Arnab, Shyamal Buch, Shen Yan, Austin Myers, Xuehan Xiong, Arsha Nagrai, and Cordelia Schmid. Streaming dense video captioning. In *CVPR*, 2024b. [1](#)
- Orr Zohar, Xiaohan Wang, Yann Dubois, Nikhil Mehta, Tong Xiao, Philippe Hansen-Estruch, Licheng Yu, Xiaofang Wang, Felix Juefei-Xu, Ning Zhang, et al. Apollo: An exploration of video understanding in large multimodal models. *arXiv:2412.10360*, 2024. [3](#)

Appendix

A INFERENCE TIME

In Table 6, we evaluated the speed of our model on a single NVIDIA A100 GPU. We employed LLaVA-Video-7B Zhang et al. (2024e) as our answering LLM, implemented a 32-frame sampling strategy from 512 input frames in total, and generated 27 text tokens. Additionally, we leveraged KV Cache and Flash Attention Dao et al. (2022); Dao (2024) to enhance inference efficiency.

Our detailed analysis of computational costs reveals that processing each video sample requires a total of 6.42 seconds, with the Vision Encoder (2.92 seconds) and LLM (2.89 seconds) dominating the time consumption. These two components collectively consume 90% of the total processing time, indicating the direction for future system optimization. In contrast, our VideoITG module demonstrates remarkable efficiency, requiring only 0.61 seconds to scan 512 frames—a speed that surpasses human visual recognition and thinking capabilities.

Table 6: Computation cost of the model.

Vision Encoder	VideoITG	LLM	Overall
2.92s	0.61s	2.89s	6.42s

B DATASET DETAILS

B.1 PROMPT TEMPLATE

Our Question-guided Clip Retrieval process utilizes a carefully designed prompt template (shown in Table 8) that instructs the LLM to analyze chronologically ordered clip-level descriptions and identify the minimal set of clips necessary to answer a given question. The prompt template consists of three main components:

- **Task Description:** Defines the LLM’s role as an expert in analyzing video clip descriptions and establishes the goal of selecting clips that cover both question and answer content.
- **Guidelines:** Provides detailed instructions for clip selection, including handling time-related expressions, determining if a single or multiple clips are needed, addressing questions about object existence or movement, and avoiding unnecessary clips.
- **Output Format:** Specifies the required JSON structure for responses, ensuring consistent formatting with explanation and clip number fields.

This template enables the LLM to perform chain-of-thought reasoning when selecting relevant clips. The model analyzes keywords from questions, identifies temporal relationships (e.g., “before,” “after”), and provides explicit rationales for its selections. For cases where no relevant clips exist, the model returns “None” to reduce annotation noise.

We implement this process using GPT-4o-mini OpenAI (2024), which is sufficient for accurate clip selection while reducing annotation costs by over 10 times compared to larger models. The selected clips are then converted to event boundaries defined by timestamps based on frame indices for the final temporal grounding annotations.

B.2 HUMAN-IN-THE-LOOP VERIFICATION

Ensuring the quality of automatically annotated datasets is critical for the reliability and effectiveness of downstream video understanding models. In this work, we implement a comprehensive quality control protocol for the VideoITG-40K dataset.

Table 7: Dataset quality (IoU). We evaluate the performance in both multiple-choice (MC) and open-ended (OE) questions.

Method	Semantic-MC	Semantic & Motion-OE	Semantic-MC	Semantic & Motion-OE
Qwen2.5-VL-32B	0.31	0.36	0.27	0.37
GPT4o	0.24	0.30	0.26	0.27
Ours	0.79	0.74	0.72	0.69

Table 8: Prompt Template: An expert system for temporal localization in video segments. The system analyzes video segment descriptions to determine the minimal and necessary combination of segments required to answer questions.

Task:

You are an expert in analyzing video clip descriptions. Your task is to select which clip or combination of clips is necessary to answer the given question, ensuring the selected clips effectively cover the content of both the question and the answer.

Guidelines:

- Carefully read the descriptions to determine which clip(s) provide relevant content for the question and the answer.
- Clip descriptions are in chronological order. Use clip number to locate clips based on time-related expressions (e.g., "at the beginning of the video" suggests a smaller clip number, while "at the end of the video" suggests a larger one).
- First, determine if one clip can answer the question or if multiple clips are needed. Then, return a list containing the selected clip(s) and an explanation.
- If the question asks about the existence/movement of an object or event. The object/action/movement may not exist, meaning you can't find the answer in the description, but the question might still provide some clues. You need to find the sentence closest to those clues.
- If asked about the whole video description or overall atmosphere, you should return all clip numbers.
- If multiple clips provide similar descriptions of the content and any of them can be used to answer the question, return all corresponding clips.
- If there are no clues in all descriptions and cannot answer the question, return "None".
- **Important:** Avoid including unnecessary clips.

Output Format:

1. Your output should be formed in a JSON file.
2. Only return the Python dictionary string.

For example:

```
{ "explanation": "...", "clip_num": "One clip: [Clip-2]" }
{ "explanation": "...", "clip_num": "Multiple clips: [Clip-1, Clip-7, Clip-8]" }
{ "explanation": "...", "clip_num": "None." }
```

Our pipeline begins with diverse sampling: we select a representative subset of the dataset, covering a wide range of instructions and video scenarios. For this subset, we conduct human verification, where expert annotators review the automatically generated annotations to assess their accuracy and relevance. This process allows us to identify and correct potential errors, and to further calibrate our annotation pipeline for improved consistency and quality.

As shown in Table 7, we compare our pipeline with baselines where advanced models such as Qwen2.5VL and GPT-4o are directly prompted to answer the temporal boundaries of relevant events.

Table 9: Prompt template for identifying motion-related questions in video QA tasks. The template instructs the system to analyze each question-answer pair and determine whether the question pertains to absolute or relative speed, responding with “Yes” or “No” accordingly. Example cases are provided for clarification.

<p>Task: Analyze the given QA pair to determine if the question is related to speed. Specifically, check if it involves either absolute speed (the speed of a specific object) or relative speed (comparing the speed of different objects). Provide an output of “Yes” if the question pertains to speed, and “No” otherwise.</p> <p>Important: Respond with “Yes” or “No” only.</p>
<p>Example: Question 1: Which is faster, the white car or the bicycle? Options: A. The bicycle. B. The white car. C. Both are at the same speed. D. None of the above. Answer 1: B. The white car. Output: Yes. Question 2: What color is the cat ?Options: A. black B. white C. orange D. gray Answer 2: C. orange Output: No.</p>

Table 10: Prompt template for identifying semantic-related questions in video QA tasks. The template instructs the system to analyze each question-answer pair and determine whether the question pertains to absolute or relative speed, responding with “Yes” or “No” accordingly. Example cases are provided for clarification.

<p>Task: Analyze the given QA pair to determine if the question inquires about the existence of an object or action. If it does, and the answer is “No” (indicating non-existence), output “Yes.” If the question is not about existence, or the answer is “Yes” (indicating existence), output “No.”</p> <p>Important: Respond with “Yes” or “No” only.</p>
<p>Example: Question 1: After going through the bag, does the person meticulously clean the area around the sink? Answer 1: No, the person does not clean the area around the sink after going through the bag. The video primarily focuses on the action of the person with the bag and items, not on cleaning activities. Output: Yes. Question 2: Is there a cat sitting on the windowsill in the video? Answer 2: Yes, there is a cat sitting on the windowsill throughout the video. Output: No.</p>

These direct approaches result in significantly lower performance, highlighting the advantage of our multi-step, instruction-guided annotation strategy.

C VISUALIZATION

In Fig. 5 and Fig. 6, we present two sets of results comparing sampling results of VideoITG with uniform sampling. Fig. 5 demonstrates a temporal reasoning problem, where our model accurately identifies the “workout” mentioned in the question and successfully locates the subsequent actions in the video, leading to the correct answer selection. In contrast, the uniform sampling strategy failed to capture these crucial frames. Fig. 6 illustrates a non-existence question scenario where our model

Algorithm 1 Keyframe Extraction via Bidirectional CLIP Similarity

Require: Video frame sequence `frames`, similarity thresholds t_1 (scene change) and t_2 (diversity)

Ensure: Selected keyframe indices `sel`

```
1: Initialize sel with the first frame index: sel  $\leftarrow$  {0}
2: Extract CLIP feature for the first frame: prev  $\leftarrow$  clip(frames[0])
3: for each frame  $c$  in frames[1:] with index  $i$  do
4:   curr  $\leftarrow$  clip(c)
5:    $s \leftarrow \text{sim}(\text{curr}, \text{prev})$ 
6:   if  $s < t_1$  then
7:     for each future frame  $f$  in frames[i+1:] do
8:       fut  $\leftarrow$  clip(f)
9:       if  $\text{sim}(\text{curr}, \text{fut}) < t_2$  then
10:        Add index  $i$  to sel
11:        prev  $\leftarrow$  curr
12:        break
13:      end if
14:    end for
15:  end if
16: end for
17: if  $\text{sim}(\text{clip}(\text{frames}[-1]), \text{prev}) < t_1$  then
18:   Add last frame index to sel
19: end if
20: return sel
```

effectively identifies all IMAX movies present in the given options, enabling it to successfully filter out and determine the correct answer.

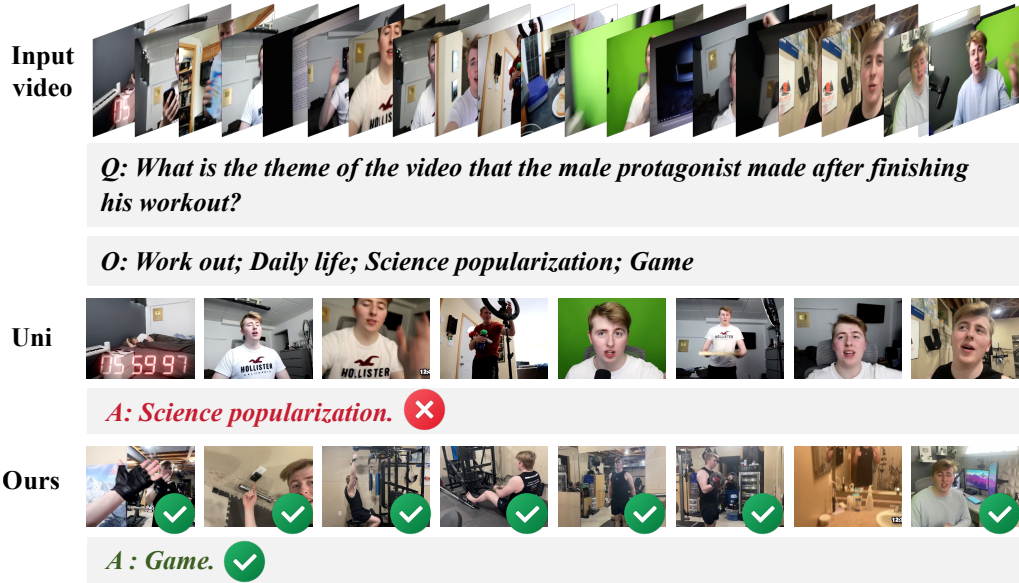


Figure 5: Example-1 shows how different sampling strategies impact video understanding. We mark the identified key frames that directly answer the question with green check-marks.

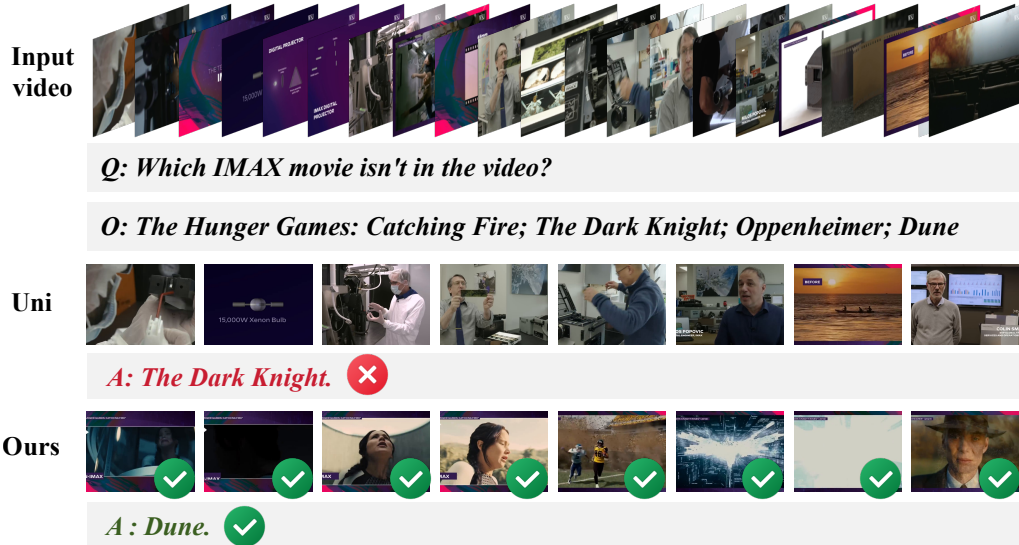


Figure 6: Example-2 shows how different sampling strategies impact video understanding. We mark the identified key frames that directly answer the question with green check-marks.